



Metadata a metadatové standardsy užívané v knihovnách

Pavla Švástová

Moravská zemská knihovna

9.12.2010 Univerzita Pardubice

1. **Metadata a metadatová schémata – teorie**
2. **Knihovnické metadatové standardy**
3. **Výroba metadat v praxi**
4. **Digitální knihovny**
5. **Využití v praxi – zajímavé projekty**

/ data o datech ???

- tuto definici lze použít v případě, že chcete člověka ještě více zmást :)
- vyjadřuje teoretickou podstatu, ALE:
- nepomůže pochopit, co to vlastně metadata jsou a k čemu slouží
- jak se liší data a metadata?









- / **„metadata is constructed, constructive and actionable“** (Understanding the Semantic Web: Bibliographic Data and Metadata, Karen Coyle
<http://alatechsource.metapress.com/content/p3022442071g7655>)
- / Constructed – metadata jsou uměle vytvořená, nenacházejí se v přírodě, jsou nadstavbou nad něčím jiným, lidskou invencí
- / Constructive – metadata jsou vytvářena cíleně, aby vyřešila nějaký problém
- / Actionable – cílem metadat je, aby byla smysluplně využita

/ aspekty metadat:

- sémantika – co chceme popsat
- syntaxe – jak to chceme zapsat
- struktura – jak vyjádřit případné vztahy

/ typ metadat:

- popisná (deskriptivní, bibliografická)
- strukturální
- technická
- administrativní

/ další možná třídění:

- jednoduchá x složitá

- / knihovnictví je od počátků založeno na vytváření a využívání metadat
- / knihovní metadata řeší problém, jak se vyznat v rozsáhlejších sbírkách dat
 - přírůstkové katalogy – seznam majetku knihovny
 - lískové katalogy (rozdělení dle autorů, názvů, témat, jazyků...)
- / vznik schémat – jak popsat knihu? co zažívá knihovník při popisu knížky?
- / http://aleph.mzk.cz:80/F?func=direct&doc_number=000059616&local_base=MZK01&format=999

- / soubor prvků určitého formátu, díky němuž lze metadata zaznamenat a vytvořit jejich strukturu
- / důležité je využívat mezinárodní standardy pro vyjádření a výměnu a vytvořit metadatové schéma tak, aby bylo využitelné pokud možno jednoznačně
- / v syntaxi dnes jednoznačně vede XML
 - specifikace formátu
 - DTD
 - jmenné tagy

- / původně vznikaly fyzicky vyjádřená metadata o fyzických dokumentech (kniha byla zkatalogizována a její metadata zapsána na katalogizační lístek)
- / nástup počítačového zápisu a zpracování fyzických knihovních jednotek – MARC
- / zdigitalizované dokumenty a born digital dokumenty – nové potřeby možností popisu – DC, MODS, METS

- / popisná = bibliografická = deskriptivní = ta klasická, co se využívají odjakživa pro popis fyzických dokumentů uložených v knihovně
- / popis dokumentů:
 - jmenný popis – autor, název, rok vydání, nakladatel, ...
 - věcný popis – pokus o popis obsahu či tématu dokumentu
 - určuje, kam je potom kniha v knihovně zařazena (předmětová hesla, MDT, klíčová slova, atd.)
- / rozdíly mezi popisem fyzických a digitálních dokumentů nejsou tak dramatické

Co nám řekne katalogizační lístek

Číslo mezinárodního desetinného třídění (MDT)

podle tohoto čísla je řazena **naučná literatura** na volném výběru

Jméno autora

podle prvních dvou písmen příjmení je řazena **beletrie** na volném výběru

Klíčová slova

92Josef II.
MAGENSCHAB, Hans
Josef II. : revolucionář z boží milosti / Hans Magenschab ; [z německého originálu ... přeložil Milan Tvrdlík]. -- 1. vyd. -- Praha : Brána : Knižní klub, 1999. -- 241 s., [8] s.obr.příl. ; 21 cm. -- Originál: Josef II. , Revolutionär von Gottes Gnaden.
Josef II.-panovníci-císařové-Rakousko-Uhersko-Čechy dějiny-Habsburkové-život-činnost-životopisy
Netradičně pojatý portrét Josefa II., jednoho z nejvýznamnějších panovníků v evropských dějinách.
80-7243-044-0 : 229.00 Kč

Signatura

podle tohoto čísla je řazena kniha , je-li ve skladu

Stručná anotace

Cena knihy

/ i ten je standardizovaný! (<http://www.knihomol.wz.cz/listek.php>)

/ projekty na retrokonverzi knihovních lístků - Retrokon

- / <http://www.loc.gov/marc/>
- / Machine Readable Cataloguing – strojově čitelný – umožňuje strojové zpracování a výměnu
- / je vyvíjen jako standard pro výměnu katalogizačních záznamů od poloviny 60.let Kongresovou knihovnou, další rozšíření nastalo v 70. letech
- / vzniklo mnoho tzv. národních formátů se základem MARC, později byl vytvořen mezinárodní formát MARC 21 (v ČR přechod z UNIMARC v roce 2004)
(<http://www.ikaros.cz/node/1761>)

FMT	BK
LDR	-----nam-a22-----a-4500
001	000059384
003	CZ-BrMZK
005	19960313143423.0
008	960301s1996---xr----f-----cze-d
020	a 80-7169-170-4
040	a BOA003 b cze
080	a 681.3.066 Linux 2 MRF
080	a 681.3.066 UNIX 2 MRF
1001	a Brandejs, Michal 7 jx20070925011 4 aut
24510	a UNIX - Linux : b praktický průvodce / c Michal Brandejs
24633	a Linux
250	a Vyd. 1.
260	a Praha : b Grada, c 1996
300	a 340 s. ; c 23 cm
500	a Bibliogr. - Rejstř.
6530	a Linux a UNIX
910	b TK-0258.237

- / základem zápisu jsou tři prvky:
- struktura záznamu
 - označení obsahu – třímístné číselné kódy a jejich podpole
 - obsah záznamu (definován dalšími standardy – ISBD, AACR2)

- / MARC 21 ve formátu XML
- / popisná metadata v digitálních knihovnách jsou uchovávána ve formátu XML
- / potřebný při konverzích mezi schématy

MARCXML



```
</datafield>
<datafield tag="020" ind1=" " ind2=" ">
  <subfield code="a">0152038655 :</subfield>
  <subfield code="c">$15.95</subfield>
</datafield>
<datafield tag="040" ind1=" " ind2=" ">
  <subfield code="a">DLC</subfield>
  <subfield code="c">DLC</subfield>
  <subfield code="d">DLC</subfield>
</datafield>
<datafield tag="042" ind1=" " ind2=" ">
  <subfield code="a">lcac</subfield>
</datafield>
<datafield tag="050" ind1="0" ind2="0">
  <subfield code="a">PS3537.A618</subfield>
  <subfield code="b">A88 1993</subfield>
</datafield>
<datafield tag="082" ind1="0" ind2="0">
  <subfield code="a">811/.52</subfield>
  <subfield code="2">20</subfield>
</datafield>
<datafield tag="100" ind1="1" ind2=" ">
  <subfield code="a">Sandburg, Carl,</subfield>
  <subfield code="d">1878-1967.</subfield>
</datafield>
<datafield tag="245" ind1="1" ind2="0">
  <subfield code="a">Arithmetic </subfield>
  <subfield code="c">Carl Sandburg ; illustrated as an anamorphic adventure by Ted Rand.</subfield>
</datafield>
<datafield tag="250" ind1=" " ind2=" ">
  <subfield code="a">1st ed.</subfield>
</datafield>
<datafield tag="260" ind1=" " ind2=" ">
  <subfield code="a">San Diego :</subfield>
  <subfield code="b">Harcourt Brace Jovanovich,</subfield>
  <subfield code="c">c1993.</subfield>
</datafield>
<datafield tag="300" ind1=" " ind2=" ">
  <subfield code="a">1 v. (unpaged) :</subfield>
  <subfield code="b">ill. (some col.) ;</subfield>
  <subfield code="c">26 cm.</subfield>
</datafield>
<datafield tag="500" ind1=" " ind2=" ">
  <subfield code="a">One Mylar sheet included in pocket.</subfield>
</datafield>
```

- / výhodou a zároveň nevýhodou je jeho rozsáhlost a náročnost – na jednu stranu je možné popsat dokument do nejmenšího detailu, nevýhodou je náročnost na zaškolení katalogizátorů
- / pro digitální objekty nebyl původně zamýšlen, takže není úplně vhodný
- / přechod na MARC ve formátu XML – umožňuje výměnu i jinými cestami, než jen přes knihovnický protokol Z39.50, např. OAI PMH
- / prakticky: pokud digitalizujeme objekt, který je již v knihovním katalogu popsán ve formátu MARC 21 konvertujeme jej do MODS a dogenerujeme další potřebné prvky

- / formát vznikl v roce 1995 z požadavků na minimální univerzální popis dokumentu
- / požadavky:
 - dostatečná jednoduchost
 - zároveň dostatečný rozsah, schopnost popsat zdroj detailně
 - možnost popisu prakticky jakéhokoliv zdroje
 - interoperabilita mezi dalšími formáty
- / vznikl standard, který tvoří 15 základních metadat prvků
- / rozšíření je možné na tzv. kvalifikovaný Dublin Core – základem je DC, rozšíření záleží na tvůrci standardu
- / <http://dublincore.org/>
- / http://www.ics.muni.cz/dublin_core/

/ Obsah:		/ Intelektuální vlastnictví:		/ Zdroj:	
/	Název	/	Tvůrce	/	Datum
/	Předmět a klíčová slova	/	Vydavatel	/	Formát
/	Popis	/	Příspěvatel	/	Identifikátor
/	Typ zdroje	/	Autorská práva	/	Jazyk
/	Zdroj				
/	Vztah				
/	Pokrytí				

http://www.webarchiv.cz/generator/dc_generator.php

/ Europeana Semantic Elements

/ <http://www.google.cz/url?sa=t&source=web&cd=1&sqi=2&ved=0CBoQFjAA&url=>

/ metadatové schéma založené na Dublin Core pro publikování v evropské digitální knihovně Europeana <http://www.europeana.eu/portal/>

/ standard obsahuje oproti DC několik prvků navíc (např. odkaz na jeden nebo několik náhledů na digitalizovaný dokument)

Metadata Object Description Schema (MODS)



- / schéma bylo představeno v roce 2002 jako kompromis mezi složitým popisem v MARC a nedostatečným popisem v DC
- / důraz na popis **digitálních zdrojů**, ale lze jím popsat prakticky jakýkoliv dokument
- / XML – eXtensible Markup Language
<http://www.w3.org/XML/>, <http://www.w3schools.com/xml/default.asp>
- / interoperabilita s dalšími schématy
- / využití v digitálních knihovnách
- / <http://www.loc.gov/standards/mods/>

MODS - struktura



```
</mods>
- <mods version="3.3">
  - <titleInfo>
    <title>František Kupka</title>
    - <subTitle>
      Lysistrata : [Národní galerie v Praze - Sbírka grafiky a kresby, Grafický kabinet, Šternberský palác, 14. října 2008 - 3. května 2009
    </subTitle>
  </titleInfo>
  - <titleInfo type="alternative">
    <title>Lysistrata - František Kupka</title>
  </titleInfo>
  - <titleInfo type="alternative">
    <title>Lysistrata</title>
  </titleInfo>
  - <name type="personal">
    <namePart>Kupka, František</namePart>
    <namePart type="date">1871-1957</namePart>
  - <role>
    <roleTerm authority="marcrelator" type="text">creator</roleTerm>
  </role>
  - <role>
    <roleTerm authority="marcrelator" type="code">art</roleTerm>
  </role>
</name>
  - <name type="personal">
    <namePart>Pravdová, Anna</namePart>
  - <role>
    <roleTerm authority="marcrelator" type="code">aut</roleTerm>
  </role>
</name>
  - <name type="corporate">
    <namePart>Národní galerie (Praha, Česko)</namePart>
    <namePart>Sbírka grafiky a kresby</namePart>
  </name>
  <typeOfResource>text</typeOfResource>
  <genre authority="marcat">catalog</genre>
```

- / TEI je konsorcium, které vyvíjí standardy pro reprezentaci textů v digitální podobě
- / <http://www.tei-c.org/index.xml>
- / využívá se pro podrobný popis rukopisů, včetně přepisu plných textů
- / specifikace, která je využívána pro popis rukopisů nejen bibliografický popis, obsahuje též popis struktury a počítá s problematikou
- / TEI P5, MASTER
- / Manuscriptorium - <http://www.manuscriptorium.com/>

- / <http://www.loc.gov/standards/vracore/>
- / standard pro popis uměleckých předmětů, uměleckých objektů a audiovizuálních dokumentů (včetně architektonických staveb, soch a dalších 3D předmětů, obrazů, fotografií, filmů, atd.)
- / <http://www.archivision.com/outgoing/vracore4/examplesindex.html>

- / týkají se spíše digitalizovaných / digitálních dokumentů
- / informace např. o:
 - použitých softwarových nástrojích a jejich verzích (není .doc jako .doc)
 - hardwarových nástrojích (typ skeneru)
 - obrazových souborech (formát a jeho specifikace, rozlišení, velikost obrázku, barevnost...)
 - kódování
 - kompatibilitě

/ informace o obrazových souborech

```
<PrimaryChromaticities>
  <PrimaryChromaticities_RedX>640/1000</PrimaryChromaticities_RedX>
  <PrimaryChromaticities_RedY>330/1000</PrimaryChromaticities_RedY>
  <PrimaryChromaticities_GreenX>300/1000</PrimaryChromaticities_GreenX>
  <PrimaryChromaticities_GreenY>600/1000</PrimaryChromaticities_GreenY>
  <PrimaryChromaticities_BlueX>150/1000</PrimaryChromaticities_BlueX>
  <PrimaryChromaticities_BlueY>60/1000</PrimaryChromaticities_BlueY>
</PrimaryChromaticities>
</Energetics>
<TargetData>
  <TargetType>0</TargetType>
  <TargetID>
    <TargetIDManufacturer>Eastman Kodak</TargetIDManufacturer>
    <TargetIDName>ColorChecker</TargetIDName>
    <TargetIDNo>Version2</TargetIDNo>
    <TargetIDMedia>Ektachrome Transparency</TargetIDMedia>
  </TargetID>
  <ImageData>00001.tif</ImageData>
  <PerformanceData>http://path/to/file</PerformanceData>
  <Profiles>http://path/to/file</Profiles>
</TargetData>
</ImagingPerformanceAssessment>
<ChangeHistory>
  <ImageProcessing>
    <DateTimeProcessed>2001-10-06T00:00:00.001</DateTimeProcessed>
    <SourceData>http://path/to/file</SourceData>
    <ProcessingAgency>JJT, Inc.</ProcessingAgency>
    <ProcessingSoftware>
      <ProcessingSoftwareName>Adobe Photoshop</ProcessingSoftwareName>
      <ProcessingSoftwareVersion>version 5.5</ProcessingSoftwareVersion>
    </ProcessingSoftware>
  </ImageProcessing>
</ChangeHistory>
```

- / <http://www.loc.gov/standards/alto/>
- / standard pro ukládání informací o vzhledu a obsahu digitalizovaného dokumentu ve formátu XML
- / „nadstavba“ nad OCR, které v tomto případě není pouze .txt
- / určuje pozici slov, řádků i odstavců na stránce s přesností na desetinu milimetru
- / zaznamenává typ a velikost písma

ALTO - příklad



```
<String ID="P1_ST00127" HPOS="1378" VPOS="1901" WIDTH="21" HEIGHT="30" CONTENT="1" WC="0.20" CC="77074"/>
<SP ID="P1_SP00099" HPOS="1399" VPOS="2018" WIDTH="31"/>
<String ID="P1_ST00128" HPOS="1430" VPOS="1950" WIDTH="146" HEIGHT="45" CONTENT="Huso" SUBS_TYPE="HypPart1"
SUBS_CONTENT="Huso.Mik" WC="0.44" CC="77074"/>
<HYP CONTENT="-"/>
</TextLine>
- <TextLine ID="P1_TL00028" HPOS="172" VPOS="2023" WIDTH="158" HEIGHT="65">
  <String ID="P1_ST00129" HPOS="172" VPOS="2023" WIDTH="158" HEIGHT="65" CONTENT=".Mik" SUBS_TYPE="HypPart2"
  SUBS_CONTENT="Huso.Mik" WC="0.47" CC="4177"/>
</TextLine>
- <TextLine ID="P1_TL00029" HPOS="496" VPOS="2112" WIDTH="645" HEIGHT="71">
  <String ID="P1_ST00130" HPOS="496" VPOS="2114" WIDTH="292" HEIGHT="69" CONTENT="Okresy" WC="0.41" CC="094977"/>
  <SP ID="P1_SP00100" HPOS="788" VPOS="2183" WIDTH="48"/>
  <String ID="P1_ST00131" HPOS="836" VPOS="2112" WIDTH="305" HEIGHT="51" CONTENT="aoudní:" WC="0.33" CC="88838"/>
</TextLine>
- <TextLine ID="P1_TL00030" HPOS="162" VPOS="2211" WIDTH="1415" HEIGHT="66">
  <String ID="P1_ST00132" HPOS="162" VPOS="2216" WIDTH="387" HEIGHT="61" CONTENT="Jihlavský;" WC="0.31"
  CC="6487887089"/>
  <SP ID="P1_SP00101" HPOS="549" VPOS="2277" WIDTH="74"/>
  <String ID="P1_ST00133" HPOS="623" VPOS="2215" WIDTH="332" HEIGHT="48" CONTENT="teltésk*." WC="0.47"
  CC="772777070"/>
  <SP ID="P1_SP00102" HPOS="955" VPOS="2277" WIDTH="80"/>
  <String ID="P1_ST00134" HPOS="1035" VPOS="2211" WIDTH="370" HEIGHT="61" CONTENT="třabiéský," WC="0.38"
  CC="8718488086"/>
  <SP ID="P1_SP00103" HPOS="1405" VPOS="2277" WIDTH="80"/>
  <String ID="P1_ST00135" HPOS="1485" VPOS="2224" WIDTH="92" HEIGHT="34" CONTENT="na" SUBS_TYPE="HypPart1"
  SUBS_CONTENT="naméifskt" WC="0.14" CC="896"/>
  <HYP CONTENT="-"/>
</TextLine>
- <TextLine ID="P1_TL00031" HPOS="163" VPOS="2290" WIDTH="1415" HEIGHT="61">
  <String ID="P1_ST00136" HPOS="163" VPOS="2298" WIDTH="282" HEIGHT="45" CONTENT="méifskt" SUBS_TYPE="HypPart1"
  SUBS_CONTENT="naméifskt" WC="0.11" CC="8988988"/>
  <SP ID="P1_SP00104" HPOS="445" VPOS="2343" WIDTH="111"/>
  <String ID="P1_ST00137" HPOS="556" VPOS="2296" WIDTH="259" HEIGHT="46" CONTENT="roto▼ic" WC="0.42" CC="53778"/>
```

- / možnost vygenerování obrázku pouze z xml souboru (záleží na kvalitě OCR, jak to bude potom čitelné :)
- / rozdělení na články
- / odlišení obrázků, textu a bílého místa na stránce
- / opravy OCR - <http://trove.nla.gov.au/newspaper>
- / projekt Impact
- / další?

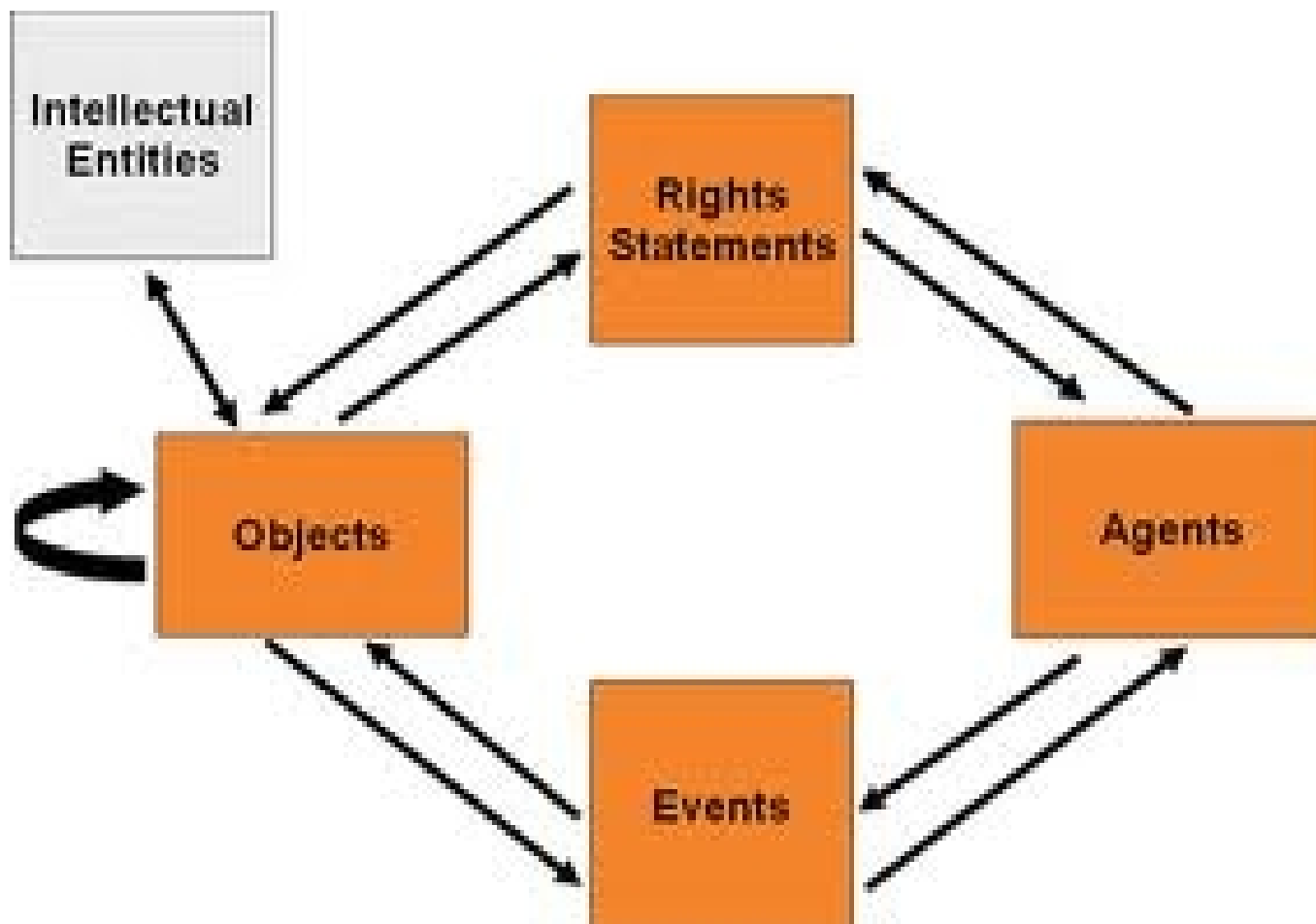
- / autorsko-právní otázky – kdo vlastní práva k dokumentu, kdy bude volně přístupný, za jakých podmínek apod.
- / prezervační funkce
- / co se s dokumentem děje za celou jeho kariéru od vzniku až po zpřístupnění – konverze do nových verzí či jiných formátů

PREMIS – The Preservation Metadata Implementation Strategies



- / <http://www.loc.gov/standards/premis/>
- / formát vznikl za účelem ochrany a dlouhodobého uchování dokumentů

- / <http://www.rinascimento-digitale.it/eventi/premis/premis-tutorial/Premis-pt1>
- / intelektuální entita (kniha, fotka, webová stránka, mapa, databáze...)
- / objekt (kapitola knihy v pdf, fotka v tiff, html soubor...)
- / událost (validace souboru, transformace, konverze, migrace...)
- / agent (autor, instituce, program...)
- / autorská práva (copyright, licence, další...)



- / digitální objekty tvoří souhrn dat a metadat, které je potřeba logicky provázat
- / př. naskenovanou a metadaty opatřenou knihu tvoří:
 - soubor obrázků v různé kvalitě a různých formátech
 - xml soubor bibliografických metadat ve formátu MODS
 - xml soubor technických metadat ve formátu MIX
 - xml soubor administrativních metadat ve formátu PREMIS
 - xml soubor OCR ve formátu ALTO

- / starší český formát metadat pro zpřístupnění v digitální knihovně Kramerius (např. <http://kramerius.mzk.cz>)
- / umožňuje popsat v podstatě jen monografii a periodikum – vytvořen účelově
- / obsahuje všechny typy metadat – popisná, strukturální i elementární technická a administrativní – jeden soubor ve formátu XML
- / dodnes se používá, ale v současnosti probíhá v NK a MZK konverze do formátů využívaných v nové verzi Krameria 4 (<http://krameriusdemo.mzk.cz>)

METS – Metadata Encoding and Transmission Standard



- / <http://www.loc.gov/standards/mets/>
- / kontejnerový formát pro uchování a výměnu metadat v digitálních knihovnách založený na XML
- / určuje strukturu dokumentu
- / provazuje všechny typy metadat v jeden složený objekt

- / může obsahovat 5 sekcí:
 - metsHdr - hlavička
 - dmdSec - popisná metadata
 - admSec - administrativní metadata
 - fileSec – objekty, soubory
 - structMap – strukturální mapa
 - behaviourSec – chování dle obsahu
- / nemusí být použity všechny sekce
- / každá sekce může obsahovat další typy metadat nebo na ně odkazovat

- / většina metadat je generovaná automaticky – software je vyrábí automaticky dle nastavení v požadovaném formátu
- / některá je třeba ručně dodělat – popřípadě opravit špatně vygenerovaná, zkonvertovaná apod.
- / <https://meditor.mzk.cz/>

- / úložiště digitálních dokumentů
- / potřebuje pro svůj provoz všechny výše zmíněné typy metadat
- / umožňuje vyhledávání a prohlížení dokumentů
- / umožňuje správu dokumentů

Příklady využití metadatových formátů v digitálních knihovnách



- / <http://kramerusdemo.mzk.cz>
- / <http://trove.nla.gov.au/newspaper>
- / <http://www.georeferencer.org/maps/domain/images>
- / <http://www.registrdigitalizace.cz>
- / <http://kramerus.mzk.cz>
- / ...

Děkuji za pozornost

Pavla Švástová
pavla.svastova@gmail.com

Moravská zemská knihovna v Brně
www.mzk.cz